



Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

Towards a unified programming model for diverse computing architectures: Experiences using PHAST

Eduardo José Gómez Hernández
Bajo la dirección de:
José Manuel García Carrasco

3 Julio 2019



Objetivos

Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

- High Performance Parallel Programming Models
- Performance Portability
- Caffe & PHAST



Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

1 Introducción

2 Revolución de los Aceleradores

3 High Performance Parallel Programming Models

4 PHAST & Caffe

5 Resultados

6 Conclusiones

7 Backup Slides



Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

- Desde el inicio, buscamos romper las barreras del rendimiento.
- La ley de Moore estableció la forma de desarrollar procesadores.
- Incrementar la frecuencia permitió aumentar el rendimiento de los circuitos.



Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

- Durante el inicio de los 2000, empezamos a encontrar un límite en la densidad energética.
- Desarrollo de un nuevo paradigma, los multi-cores.
- Aún así, habían tareas que era mejor delegar a otros dispositivos.



Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

- Al principio fue usada únicamente para gráficos.
- Paralelismo SIMT (Single Instruction Multiple Threads).
- Muy extendido como acelerador, GPGPU.



Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

- Las FPGAs eran usadas para crear prototipos de nuevo hardware.
- Ahora se empiezan a usar como hardware reconfigurable.
- Aún así, para tareas específicas los ASICs son mucho más eficientes que las FPGAs.





Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

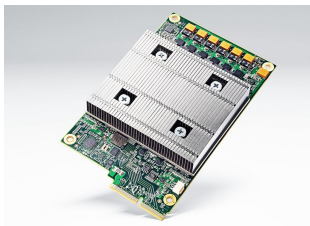
PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

- Gran tiempo sin evolucionar por falta de poder de cómputo.
- Las redes neuronales es una de las técnicas más potentes.
- Está siendo usado en muchos campos, desde medicina hasta economía.
- Desarrollados usando frameworks optimizados debido a la cantidad de cómputo requerida.
- Reciente incremento del uso de aceleradores como las TPUs o los Tensor Cores, e incluso hardware más específico.





Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

- Los clústers tradicionales estaban sufriendo graves problemas para escalar ciertas aplicaciones.
- Un clúster heterogéneo permite una mejor escalabilidad, pero conlleva varios problemas.



Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

- Uno de los problemas de los clúster heterogéneo es usarlo de manera eficiente.
- Existe una gran variedad de modelos para usar un clúster heterogéneo de manera eficiente.



Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

- De 50 modelos revisados, hemos clasificado 23.
- Hemos basado la clasificación en 4 factores.
 - Plataforma
 - Tipo de modelo
 - Paradigma de programación
 - Paradigma de memoria



Clasificación

Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel
Programming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

CPU	Programming Paradigm									Memory Paradigm								
	Language	Kernel	Task	Directive	C++ Template	Skeleton Lift	Wrapper	Threads	Hierarchical	Shared	Implicit	Explicit			PGAS	Distributed		
		OpenCL							OpenCL	Cilk	Lift							
Type	Extension		Cilk	OpenMP	OpenAcc				OpenMP	Cilk	AllScale	OmpSs						
	Library	SyCL	Kokkos	TBB	MPI	Phast	Dash	SkePU	EasyCL	Pthreads	SkePU	Phast	Pacxxv2	EasyCL	Thrust	Dash	SyCL	TBB
		Pacxxv2	MultiController	StarPU			Thrust	Boost Compute	SkeCL				Boost Compute	MultiController	StarPU	Kokkos	MPI	

GPU	Programming Model									Memory Model							
	Language	Kernel	Task	Directive	C++ Template	Skeleton Lift	Wrapper	Threads	Hierarchical	Shared	Implicit	Explicit			PGAS	Distributed	
		OpenCL							OpenCL		Lift						
Type	Extension			OpenMP	OpenAcc				OpenMP			OmpSs					
	Library	SyCL	Kokkos	StarPU		Phast	Dash	SkePU	EasyCL		SkePU	Phast	Pacxxv2	EasyCL	Thrust	Dash	SyCL
		Pacxxv2	MultiController				Thrust	Boost Compute	SkeCL				Boost Compute	MultiController	StarPU	Kokkos	

FPGA	Programming Model									Memory Model						
	Language	Kernel	Task	Directive	C++ Template	Skeleton	Wrapper	Threads	Hierarchical	Shared	Implicit	Explicit			PGAS	Distributed
		OpenCL							OpenCL							
Type	Extension											OmpSs				
	Library	SyCL						SkeCL	EasyCL					EasyCL		SyCL
														SkeCL		



Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

- Modelos capaz de ser ejecutados en múltiples dispositivos usando el mismo código fuente.
- Ejemplos:
 - OpenMP (Algunas implementaciones)
 - OpenAcc
 - OpenCL & SYCL
 - SkePU
 - Kokkos
 - PHAST



Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

- Modelos capaz de ser ejecutados en múltiples dispositivos usando el mismo código fuente.
- Ejemplos:
 - OpenMP (Algunas implementaciones)
 - OpenAcc
 - OpenCL & SYCL
 - SkePU
 - Kokkos
 - PHAST



Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

- Librería de C++ con estructura similar a la librería algorítmica estándar de C++.
- Capaz de generar código para GPUs de Nvidia, y para multi- y many-core.
- Uso de `functors` para modificar los contenedores de datos.
- Define un modelo de programación agnóstico a la plataforma, pero permite el uso de porciones de código específico.

PHAST



Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

- Uno de los primeros frameworks de Deep Learning, desarrollado por Berkeley.
- Ha sido uno de los más usados hasta que terminó su desarrollo.
- Sigue siendo usado para investigación, gracias a su facilidad de ser modificado.



Roadmap

Introducción

Revolución
de los
Aceleradores

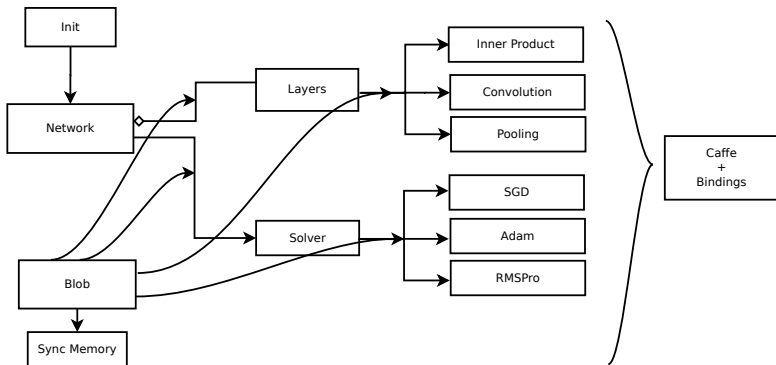
High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides





Roadmap

Introducción

Revolución
de los
Aceleradores

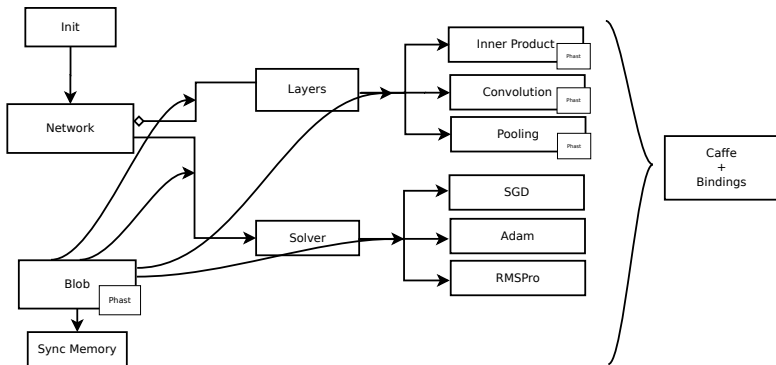
High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides





InnerProduct Caffe CPU

Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

```
1  template <typename Dtype>
2  void InnerProductLayer<Dtype>::Forward_cpu(const vector<Blob<Dtype>*>& bottom,
3         const vector<Blob<Dtype>*>& top) {
4
5     const Dtype* bottom_data = bottom[0]->cpu_data();
6     Dtype* top_data = top[0]->mutable_cpu_data();
7     const Dtype* weight = this->blobs_[0]->cpu_data();
8
9     caffe_cpu_gemm<Dtype>(CblasNoTrans, transpose_ ? CblasNoTrans : CblasTrans,
10      M_, N_, K_, (Dtype)1.,
11      bottom_data, weight, (Dtype)0., top_data);
12
13     if (bias_term_) {
14         caffe_cpu_gemm<Dtype>(CblasNoTrans, CblasNoTrans, M_, N_, 1, (Dtype)1.,
15          bias_multiplier_cpu_data(),
16          this->blobs_[1]->cpu_data(), (Dtype)1., top_data);
17     }
18 }
```



InnerProduct PHAST

Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

```
1  template <typename T, unsigned int policy = phast::get_default_policy(>
2  struct matrixPlusVectorRows : phast::functor::func_vec<T, policy> {
3      _PHAST_METHOD matrixPlusVectorRows() {}
4
5      _PHAST_METHOD void operator()(phast::functor::vector<T>& row) {
6          for(auto r = row.begin(), i = vec.begin(); r != row.end(); ++r, ++i)
7              *r += *i;
8      }
9      phast::functor::vector<T> vec;
10 };
11
12 template <>
13 void InnerProductLayer<float>::Forward_cpu(const vector<Blob<float>*>& bottom,
14 const vector<Blob<float>*>& top) {
15     phast::matrix<float> matA = bottom[0]->getDataAsMatrix(M_, K_, false);
16     phast::matrix<float> matB = this->blobs_[0]->getDataAsMatrix(K_, N_, !transpose_);
17     phast::matrix<float> matC = top[0]->getDataAsMatrix(M_, N_, false);
18
19     phast::dot_product(matA, matB, matC);
20
21     if (bias_term_) {
22         matrixPlusVectorRows<float> matrixPlusVectorRows;
23         matrixPlusVectorRows.vec.link(this->blobs_[1]->getDataAsVector(N_));
24         phast::for_each(matC.begin_i(), matC.end_i(), matrixPlusVectorRows);
25     }
26     if(!transpose_) matB.transpose();
27 }
```

Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

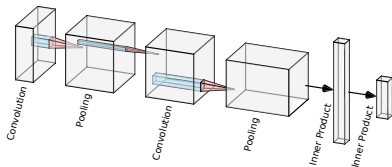
PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

- Hemos usado dos redes neuronales basadas en LeNet para comprobar el funcionamiento de nuestras modificaciones.
- Hemos usado una máquina del clúster de GACOP, 2 Xeon CPU E5-2603 v3 y una GPU Nvidia GTX 1080.
- Ambas redes han sido ejecutadas correctamente en CPU y GPU.





Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

- Caffe incluye una serie de tests para comprobar cada capa de manera individual.
- A excepción de la funcionalidad no implementada, todos los tests han pasado correctamente.



Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

- Debido a la limitación de tiempo los resultados que tenemos, no son definitivos.
- Los casos de prueba son demasiado pequeños.
- El código puede ser mejorado fácilmente.



Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

- Clasificación de los modelos de programación más relevantes.
- Extracción de algunos modelos de programación portables.
- PHAST es factible para aplicaciones reales.



Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

- Extender la clasificación de modelos.
- Completar y liberar la versión de Caffe con PHAST.
- Extender el trabajo a otros dispositivos.
- Mejorar el rendimiento.



Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

Towards a unified programming model for diverse computing architectures: Experiences using PHAST

Eduardo José Gómez Hernández
Bajo la dirección de:
José Manuel García Carrasco

3 Julio 2019



Introducción

Revolución
de los
Aceleradores

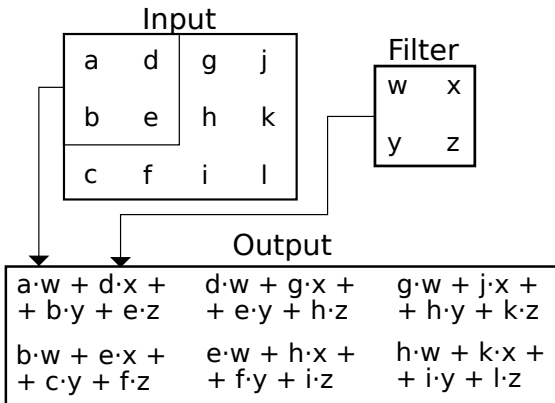
High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides



Introducción

Revolución
de los
Aceleradores

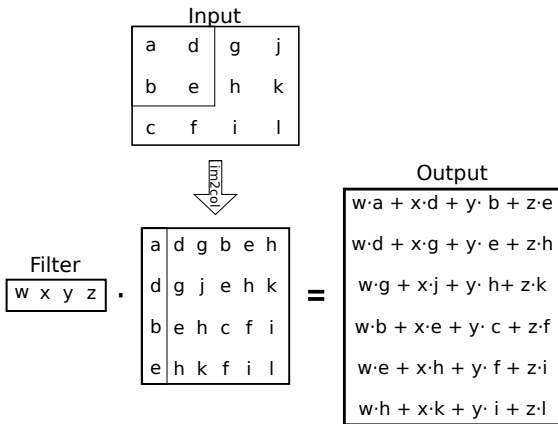
High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides





im2Col Caffe CPU

Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

```
1  const int output_h = (height + 2 * pad_h - (dilation_h * (kernel_h - 1) + 1)) /  
    stride_h + 1;  
2  const int output_w = (width + 2 * pad_w - (dilation_w * (kernel_w - 1) + 1)) /  
    stride_w + 1;  
3  const int channel_size = height * width;  
4  for (int channel = channels; channel--; data_im += channel_size) {  
5  for (int kernel_row = 0; kernel_row < kernel_h; kernel_row++)  
6  for (int kernel_col = 0; kernel_col < kernel_w; kernel_col++) {  
7  int input_row = -pad_h + kernel_row * dilation_h;  
8  for (int output_rows = output_h; output_rows > 0; output_rows--) {  
9  if (!is_a_ge_zero_and_a_lt_b(input_row, height)) {  
10     for (int output_cols = output_w; output_cols > 0; output_cols--)  
11         *(data_col++) = 0;  
12     } else {  
13         int input_col = -pad_w + kernel_col * dilation_w;  
14         for (int output_col = output_w; output_col > 0; output_col--) {  
15             if (is_a_ge_zero_and_a_lt_b(input_col, width)) {  
16                 *(data_col++) = data_im[input_row * width + input_col];  
17             } else {  
18                 *(data_col++) = 0;  
19             }  
20             input_col += stride_w;  
21         } }  
22         input_row += stride_h;  
23     } } }  
24 }
```



Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

```
1  int index = this->get_index();
2
3  const int oih = (h + 2 * ph - (dh * (kh - 1) + 1)) / sh + 1;
4  const int oiw = (w + 2 * pw - (dw * (kw - 1) + 1)) / sw + 1;
5
6  const int irow = -ph + (((index / (oih*oiw)) % (kh * kw)) / kw) * dh + ((index % (oih*
7  oiw)) / oiw) * sh;
8  const int icol = -pw + (((index / (oih*oiw)) % (kh * kw)) % kw) * dw + ((index % (oih*
9  oiw)) % oiw) * sw;
10 if (irow >= 0 && irow < h && icol >= 0 && icol < w)
    col = in.at(((index / (oih*oiw)) / (kh * kw)), irow * w + icol);
    else col = 0;
```



Blob

Introducción

Revolución
de los
Aceleradores

High
Performance
Parallel Pro-
gramming
Models

PHAST &
Caffe

Resultados

Conclusiones

Backup
Slides

